

# 知的映像メディア検索技術の動向

## Intelligent Retrieval for Video Media

美濃 導彦\*  
Michihiko Minoh

\* 京都大学工学部附属高度情報学実験施設  
Integrated Media Environment Experimental Laboratory, Kyoto University.

1995年9月18日 受理

Keywords: intelligent retrieval, video media, multimedia, pattern recognition.

### 1. はじめに

マルチメディアという言葉は、さまざまな意味で用いられるが、ここでは、音声、画像、文字データなどの異なるメディアを統合的に扱う技術の総称として用いる。

マルチメディアデータは、マルチメディアで扱われるデータの総称で、その中心は動画と静止画や動画の静止画を「画像」静止画の時系列を「映像」という用語を用いる。

人間どうしのコミュニケーションは時間と空間を共有することが前提であった。すなわち、同じ時間と空間に存在することが前提であり、コミュニケーションの始まりである。この制約のもとでは人間は五感を利用して活用できるもので、マルチメディア処理によるコミュニケーションが可能である。

しかし、見方を換えればこの制約は大きな負荷となる。すなわち、空間的に近くにいる異時代の人間とシ

\*1 マルチメディアといっても、厳密には、映像、映像は扱っている。人間の入力情報の90%程度は視覚からの情報であるといわれているので、ここでは除外して考えるが、こ

\*2 メディアのデジタル化に伴って、コンピュータでこれらの情報を扱う。的に扱うようまでの水準になってきているとい

\*3 電子情報通信学会第5回基礎研究発表会(1995年10月10日〜10月13日、京都府京都市)で発表。平成7年4月〜平成10年3月。

### 2. 映像検索の特徴

#### 2.1 映像情報空間

映像は動画と音声の両方を含んだものである。空間的には、画像を2次元、音声を1次元、合わせて3次元で、これに時間軸が加わった4次元空間である。これは、デジタル化された映像を信号として扱った考え方であり、計算機に与えられる膨大なデータの配列の次元数が4であるといっているだけである。

これにもう1次元、情報抽出の次元を加える。例1に示すように、これは、従来のパターンの認識における情報の抽象化の次元である。一番低いレベルは信号のレベル(数値の配列)で、次の段階が画像処理や音声処理(数値の配列)で、さらに抽出された特徴レベルの配列である。特徴レベルからさらに抽象化を進めようとするときマンディックギャップと呼ばれる。その際このようにあるのが人工知能の分野で扱われているシンボル記述のレベル(テキスト表現、言語表現)である。映像を見て人間がキーワードを与える場合は、人間がこの処理をするので、直接シンボルのレベルの情報が得られる。計算機による処理を行う場合は、特徴レベルとシンボルのレベルの両方を用いる「モデル」が必要である。情報抽出の次元は、マルチメディアデータに、検索のためのインデックスを付加するものである。

映像空間は、図2に示すような次元空間で表現できる。画像の2軸と、音声の1軸は、信号がそのまま

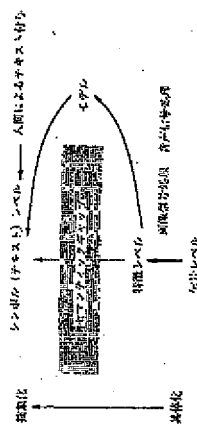


図1 映像情報空間

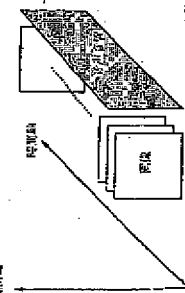


図2 映像空間

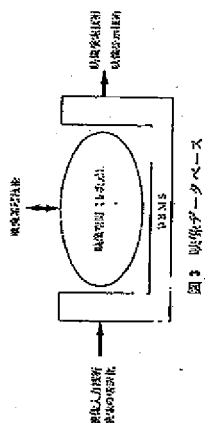


図3 映像データベース

入る軸であり、時間軸に対しては同期がとれている。現在のところ、映像は映像、音声は音声の信号として処理されている。一般的には、音声と映像を結合して処理する処理は意味づけが困難であり、現在まで行われていない。しかし、将来的には試みる価値がある問題である。

時間軸を持つメディアは、その場で処理するのと、記録したものを処理するのとで、ほとんど同じ時間が必要である。例えば、静止記録を考慮してもいい。これをデ・ブレンコードに録音したものを真剣に撮こうとすると編集と同じ時間がかかる。これが、書き出しと書き込みの時間差である。一度書き込み、書き出しが完了すれば、時間の制約から解放される。

情報抽出に対するデータは、もとのマルチメディアデータを処理して作成するか、人間がキーワードを付加することにより生成される。この種のデータの質が、知能的検索を実現するうえでのキーポイントになる。

図3に示すように、映像空間に代表されるマルチメディア空間にデータを投入する処理、空間に存在するデータを検索する処理、蓄積されているデータを検索・抽出する出力処理のすべてが、マルチメディアデータベース、特に、映像データベースには必要である。

本稿では、特に、入力側での映像の構造化、検索側の映像検索の方法、結果の提示方法に重点を置いて述べる。映像の検索技術に関しては、映像の外部記録装置への接続方法、再生時の音声と映像の同期問題、映像品質におけるリアルタイム性の維持など、システム的に興味深い問題があり、VOD技術と関連しているという点で研究が進んでいるが、ここでは議論しない。

#### 2.2 映像メディアの検索と提示

検索は映像空間に含まれる情報をいかに取り出すかという問題である。検索手法は、図4に示すように、この空間の任意の断面で表現できる。当然、クエリもマルチメディアデータを利用して行う。映像に刻し

ては以下のような検索手法が考えられる。

う形でクエリーが発せられたかにも依存する。また、検索がうまくいかなかった場合、システムがどのように修正したかをユーザに示すことも重要である[前田 95]。

### 2.3 知的映像検索

「知的な検索」とは何を意味するのであるか? クエリーが意味でも適当に検索してくれるような検索を想像する人もいれば、ある種の推論を問にはさんで検索してくれるようなシステムを想像する人もいて、思われる。本稿では、データベースへのデータの登録時にそれほど手間をかけなくても、それなりに検索できるようなメカニズムを構築して「知的な検索」と考える。この意味での知的な検索ができるシステムを実現するうえでの技術的な問題は次の3点である。

(1) 情報として与えられるマルチメディアデータから検索手法を抽出したさまざまな検索のためのインデックスを作成する技術、言い換えるならば、映像等の情報抽出・表示に対するデータを自動生成する技術、これはパターン認識の分野の技術である。この処理は、認識誤りが本質的に完全な認識は期待できないことを認識しなければならない。

(2) インデックスが認識誤りを含んでいても、大きな検索ミスを起こさないような検索手法が必要である。完全なデータベースから不完全データベースへと進んできたクエリー処理を、誤りを含まずデータベースに拡張することが必要である。

(3) データの類似度に関する度合い(類似度)の定義がユーザによって異なるので、個人に適合できるように、検索の過程において類似計算が定義できるメカニズムが必要である。また、システムが行っている検索の内容をユーザに提示する手法も、知的検索システムには不可欠である。

現在報告されている映像処理、検索システムでは、これらの点はそれほど意識されてはいないが、今後は重要になると考える。

### 3. 映像の構造化技術

「延の映像」は、図5に示すように、1枚の静止画像とは違って、作成者が選択したストーリーを伴っている。このストーリーは、いくつかのサブストーリーに分けられ、さらにサブストーリーも分割できる。このように映像は、一般的には階層構造をしており、ツリー構造データで記述できる[奥村 93]。

このツリーの葉にあたる部分は「連の意味を持つ連

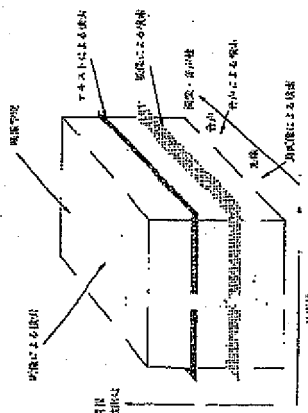


図4 検索の手法

映像による検索  
音声による検索  
映像による検索(動画像だけでなく、音声も含む)  
これらの検索手法は排他的ではなく、共存が可能である。例えば、男女がバレーを歩いている映像を検索する場合のクエリーとして、  
クエリー1: テキサスト: 「男と女が 一緒に歩いている」  
クエリー2: 画像: 男女がバレーの映像  
クエリー3: 動画像: 男女がバレーを歩いている動画像

などが考えられる。これらを同時に与えれば、システム側は有効な情報が多くなるので、良い検索結果を提示できるはずである。このように考えると、クエリーはデータのフィルタの役目を果たしている。これは画像データベース場合とまったく同様である[Paloutos 94, Mukunoki 94a]。

クエリーをどういう形でシステムに与えるかは問題である。テキストで与える場合は、マルチメディアオブジェクトとその関係を言葉で記述する。当然、人間の動作や音の特徴なども記述可能である。対象物の動きや色、形などをもとに検索する場合は、言葉で表現するよりも、データそのものをクエリーとして利用する方が便利である。この場合のクエリー処理は、与えられたメディアデータに対して「似たもの」を探し処理となる。一般的には、データにはさまざまな解像度が存在するので、逆変換するが、大問題である。次の問題は、検索された映像データは、その内容を認識するためには、それなりの時間がかかる。検索されたものが目的に合っているかどうかを判断するためにはブラウジング機能が必須になる。これは、どうい

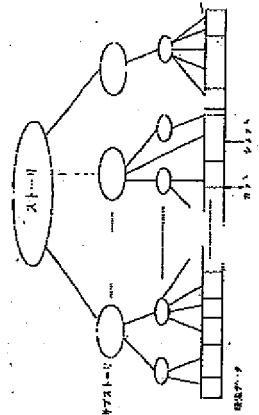


図5 映像データの構造

縮したフレームの系列であり、ショットと呼ばれる。ショットの切れ目をカットと呼ぶ。これらの単位は同じことを通った側面から表している。

映像を処理・要約して知的検索を実現するためには、基本的には、映像をショットに分割する必要がある。映像を計算機で扱う利点は、VOIなどのように同時アクセスに対するサードの多変数やほぼ映像の異なる部分の再生など、従来の記述媒体に比較しているところがあるが、やはり、意味的な処理を実現するためには意味のある最小単位であるショットを抽出する処理は不可欠である。

映像の構造化処理としては、ショットを抽出した後、それらを分類し、インデックスを付加しておく必要がある。そのためには、ショットから特徴量を抽出し、ショットを分類する処理が必要となる。

基本的には、情報抽出のデータを作成する処理は、映像を対像としたデータベースで付与されたものをショットに対して適用する処理になる。また、映像は、同時に音声や時間軸を持つており、次の数が多いので、データ量も多く、計算機で扱うのが画像ほど簡単ではない。これに伴って、検索方法も多様化する。また、結果を提示する方法も映像特有のものが必要であり、ヒューマンインテグレーションに重点を置いた研究が所望に行われている。

### 3.1 映像からのカット抽出法

人間が映像を見た場合、カットはかなりの正確に検出できる。この操作において、人間は意味的な情報と信号号的な情報を統合的に利用していると考えられる。

計算機にカット検出を行わせる場合は、信号的な情報のみに頼らなければならない。そのためには、カットとは何かを映像の信号レベルの記述に置き換えるカットのモデルを構築する必要がある。もちろん、常に、ショットのモデルを考えても差し支えない。

カットのモデルとしては以下のようなものがある

[大辻 92]

(1) カットの前後で画面の分布が時間的に不連続に変化する。

(2) きわめて短い時間に大きな変動量が発生し、直後に収束する。

(3) 変動は画面の広い範囲に及ぶ。

多くのカット抽出法は、これらのモデルの一つまたは複数のモデルを利用している。これらのモデルをどう実現するかでさまざまなカット抽出法が提案されている。例えば、(1)のモデルを利用する場合でも、画面の輝度ヒストグラムの差分を用いる方法[大辻 89, 長坂 92]もあれば、フレーム間差分に基づく方法[大辻 91]も考えられる。しかし、本質的に同じモデルに基づくので、手法的にそれほど大きな違いはない。これらの手法の比較検討については「大辻 92」を参照されたい。

これらの手法は、さまざまなパラメータを持つており、その調整によりカット抽出率、未検出率が変化する。「大辻 92」によれば、カット検出率・カット未検出率となるようにパラメータを調整すればどの手法でもカット検出率を3%以下に抑えられる。

ここであげたカットのモデルは、フェードアウトや特殊効果による徐々に変化するカットの抽出には失敗する。これに対処するために、カットの前後数フレームでの画像データの統計を利用した手法[Zhang 93]が報告されている。

たとえばどのようなモデルを作成しても、セマンティックギャップの問題があり、完全なモデルの生成は困難である。意味的な情報に基づいてカットが存在する限り、信号的な手法では検出はできない。したがって、これらのカット抽出法は、誤りを本質的に含むものとして、扱うことを考えていかなければならない。

### 3.2 ショットからの特徴抽出

何らかの形でショットが抽出できると、次の課題は抽出されたショットから検索や編集処理に役立つ特徴量を抽出することである。特徴量として現在までに利用されているのは、撮影時のカメラの動きに関する情報、撮影者や記者がつけたアノテーション、ショットに含まれる色情報などである。

撮影時のカメラの動きは、オブジェクトの位置に基づいて簡単に検出できる[大辻 92]。これはカメラの動きに対応して出現するオブジェクトの位置のバリエーションが反映していることを利用している。ただし、オブジェクトの位置の計算は検出が難しく、それほど信頼性は低くない。



## 論文 (Technical Papers)

A Constraint-Based Knowledge Compiler for Parametric Design Problems in Mechanical Engineering .....	Yasuo Nagai · Satoshi Terasaki	60
カルデシアン空間モデルに基づく知識獲得支援システム .....	矢口博之 A Knowledge Acquisition System Based on the Cartesian Space Model .....	H. Yaguchi 75
知識の分類を利用した類推による問題解決の効率化 .....	岡山将也 · 真野芳久 · 村本正生 Efficient Problem Solving by Analogical Reasoning Using Knowledge Classification .....	N. Okayama Y. Mano M. Muramoto 86
順序表現向き個体群探索分岐型遺伝的アルゴリズム o-f-ga (Order-Based-Forking GA) .....	岡井茂義 · 藤本好司 The o-fga : Order-Based Forking Genetic Algorithm .....	S. Tsutsui Y. Fujimoto 96
リアルタイム時相論理による並列ソフトウェアの 計算機支援設計に関する研究 .....	山根 智 Study on Computer Aided Design of Parallel Software by Real-Time Logic .....	S. Yamane 105
工学システムの機能モデルからの挙動の導出 .....	五福明夫 Deriving Behaviour of an Engineering System from a Functional Model .....	A. Gofuku 112
最小化の組合せ探索による帰納推論 .....	秋葉豊孝 · 佐藤素介 Searching Combinations of Fast General Generalizations for Inductive Inference .....	S. Akiba T. Sato 121
帰納的アルゴリズムに基づく巡回セールスマン問題の解法 .....	新妻清三郎 · 村田安永 · 山田和年 Inductive Algorithms for Traveling Salesman Problem .....	S. Niitsuma Y. Murata K. Yamada 130
A Query Evaluation Method for Abductive Logic Programming Based on Generalized Stable Models .....	Ken Satoh · Noboru Iwayama	137
物語のための技法と戦略に基づく物語の概念構造生成の 基本的フレームワーク .....	小方 孝 · 堀 浩一 · 大須賀範雄 A Basic Framework for Narrative Conceptual Structure Generation Based on Narrative Techniques and Strategies .....	T. Ogata K. Ijori S. Ohsguza 148

[illegible]

ISSN 0912-8085 発売元 才一出版社 定価 2400円(本体 2330円・税 70円)

介紹者

题词

英漢 釋摩